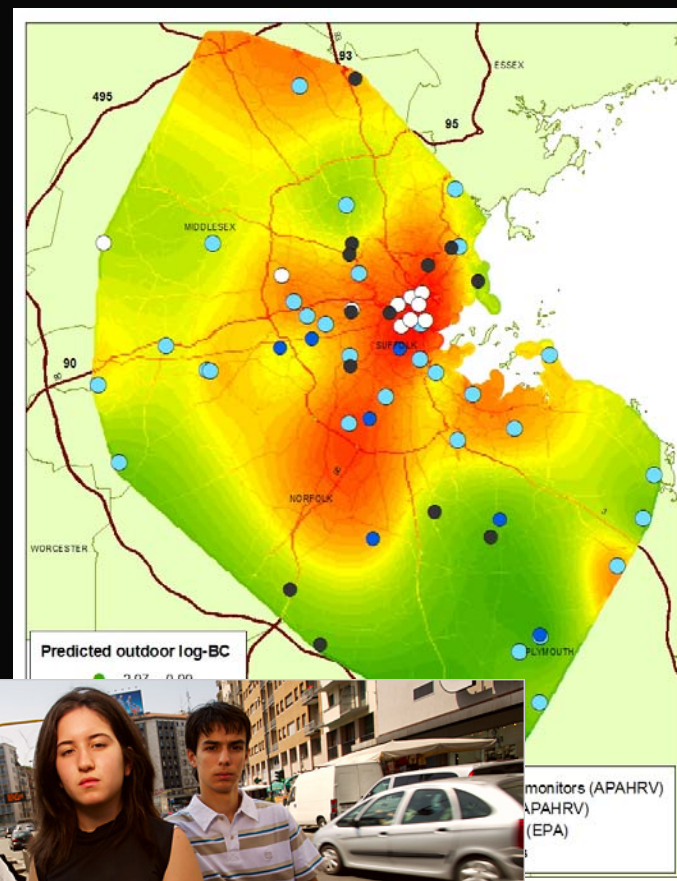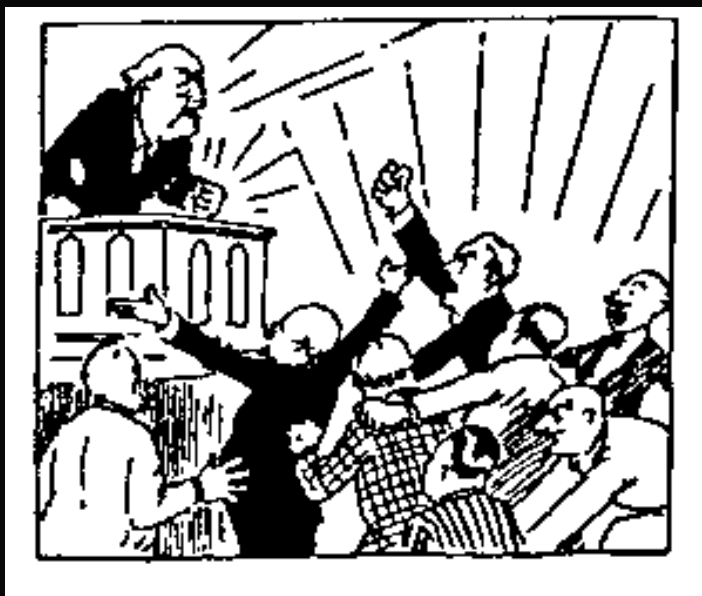# Community-based Risk Assessment – a statistician's perspective

**Louise Ryan**

**Department of Biostatistics**

**Harvard School of Public Health**



1

# Outline

❑ **Use some examples to**

– **Illustrate challenges**

– **Describe useful statistical tools and areas where more research would be helpful**

❑ **My examples**

– **Classic cancer cluster investigation**

– **Home Allergen Study**

– **Exposure assessment for various Boston based studies**

– **Mercury and IQ**
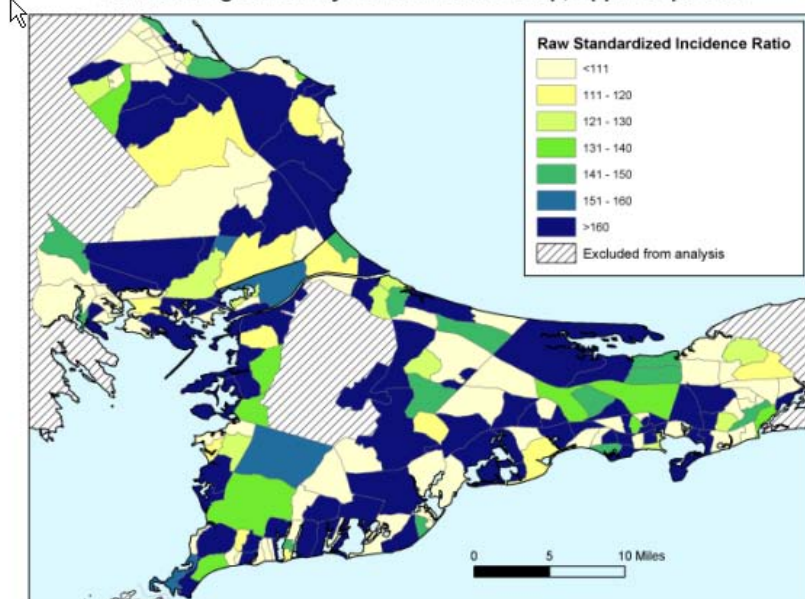
# Cancer risks on Cape Cod



❑ **Citizens near air-force base concerned about excess cancer rates reported on upper cape**

❑ **Clear evidence of multiple exposures**

❑ **Excesses small to moderate (SIRs around 120)**
  – **Power limited by total pop of ~30K**
  – **No individual exposure assessment**

3

# Cape Cod - continued



☐ **Data very noisy – smoothing no help**

☐ **Very frustrating experience for all**
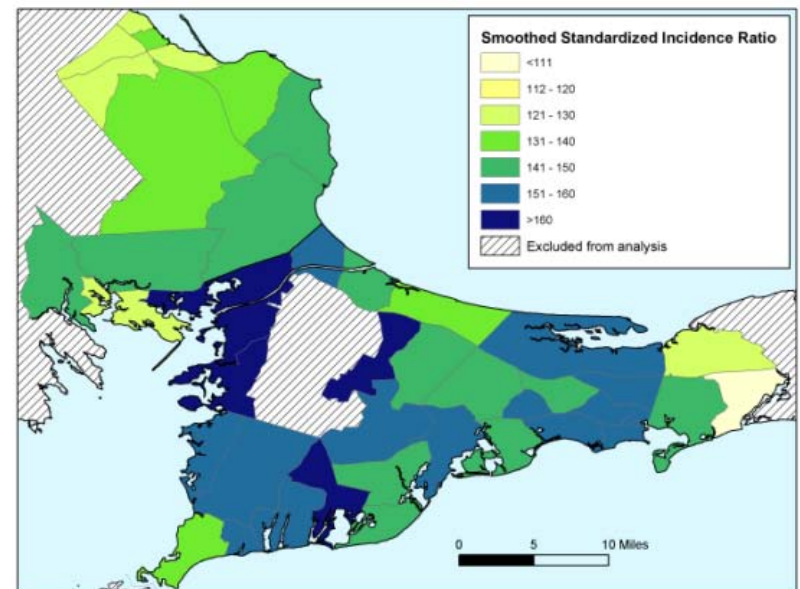
☐ **Need guidelines on what's achievable**



Female Lung Cancer by Census Block Group, Upper Cape Cod

Raw Standardized Incidence Ratio
- <111
- 111 - 120
- 121 - 130
- 131 - 140
- 141 - 150
- 151 - 160
- >160
- Excluded from analysis

0    5    10 Miles

9/14/2001          Sources: MA Dept. of Public Health, Harvard School of Public Health



Female Lung Cancer by Census Tract, Upper Cape Cod

Smoothed Standardized Incidence Ratio
- <111
- 112 - 120
- 121 - 130
- 131 - 140
- 141 - 150
- 151 - 160
- >160
- Excluded from analysis

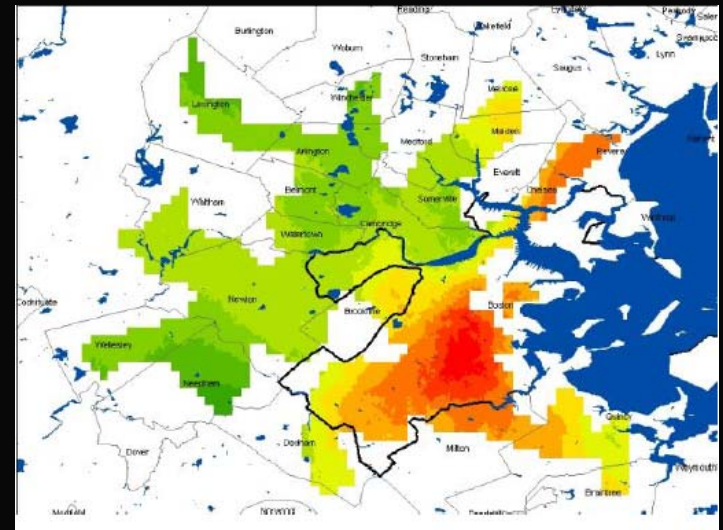0    5    10 Miles

9/13/2001          Sources: MA Dept. of Public Health, Harvard School of Public Health

# Home Allergen Study

- **Mother/child pairs recruited at birth. Followed for asthma, allergy, respiratory disease**
- **Interest in allergens, molds, adjusting for social factors**
- **Geocode study subjects and assign areal level characteristics (e.g. based on census)**

**Intriging geographical variation in maternal serum IGE. But geoadditive modeling (Kammen & Wand) suggests "hotspot" confounded with race, poverty.**

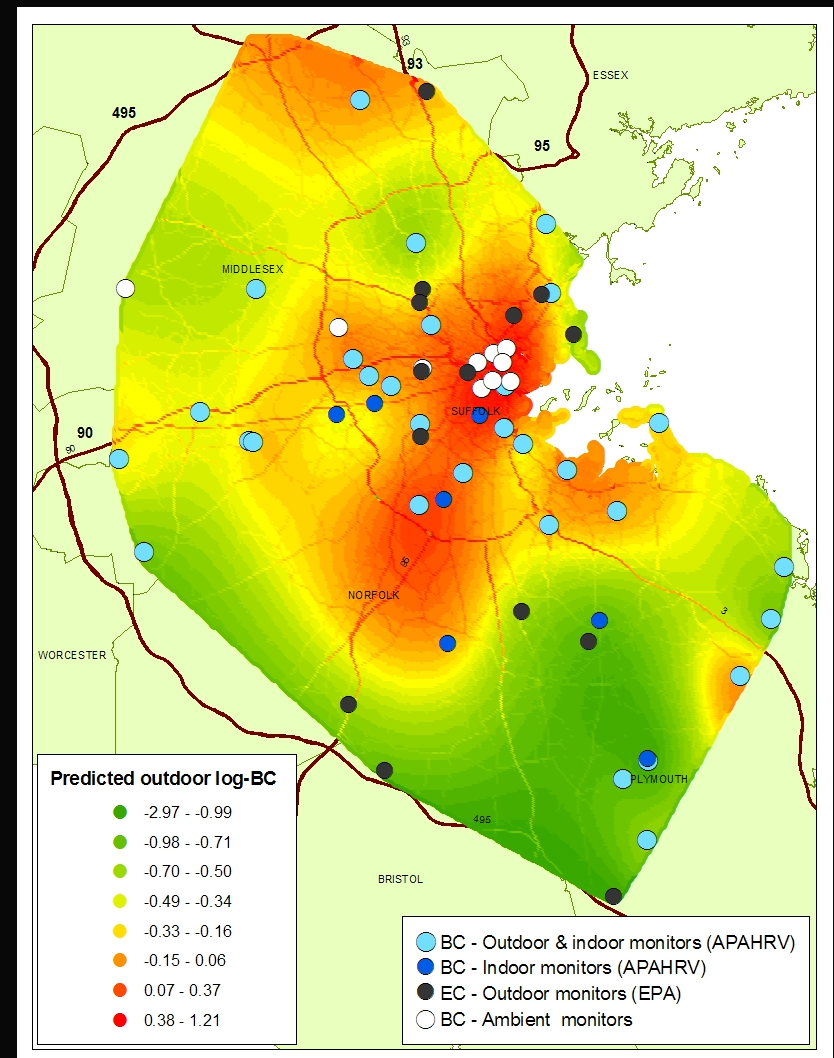$$Y = \beta_0 + \beta X_1 + g(X_2) + h(lat, lon) + \varepsilon$$

5

# Boston and New England studies of cardiovascular response to air pollution

- ❑ **Estimate exposure from**
  - – **EPA EC monitors**
  - – **Various Indoor & outdoor monitors (different studies)**
  - – **GIS-based measures (traffic density, potentially climate, land use etc)**
- ❑ **Goal – relate predicted exposures to health outcomes (heart rate variability, arythmias, birth weight), accounting for estimation error**
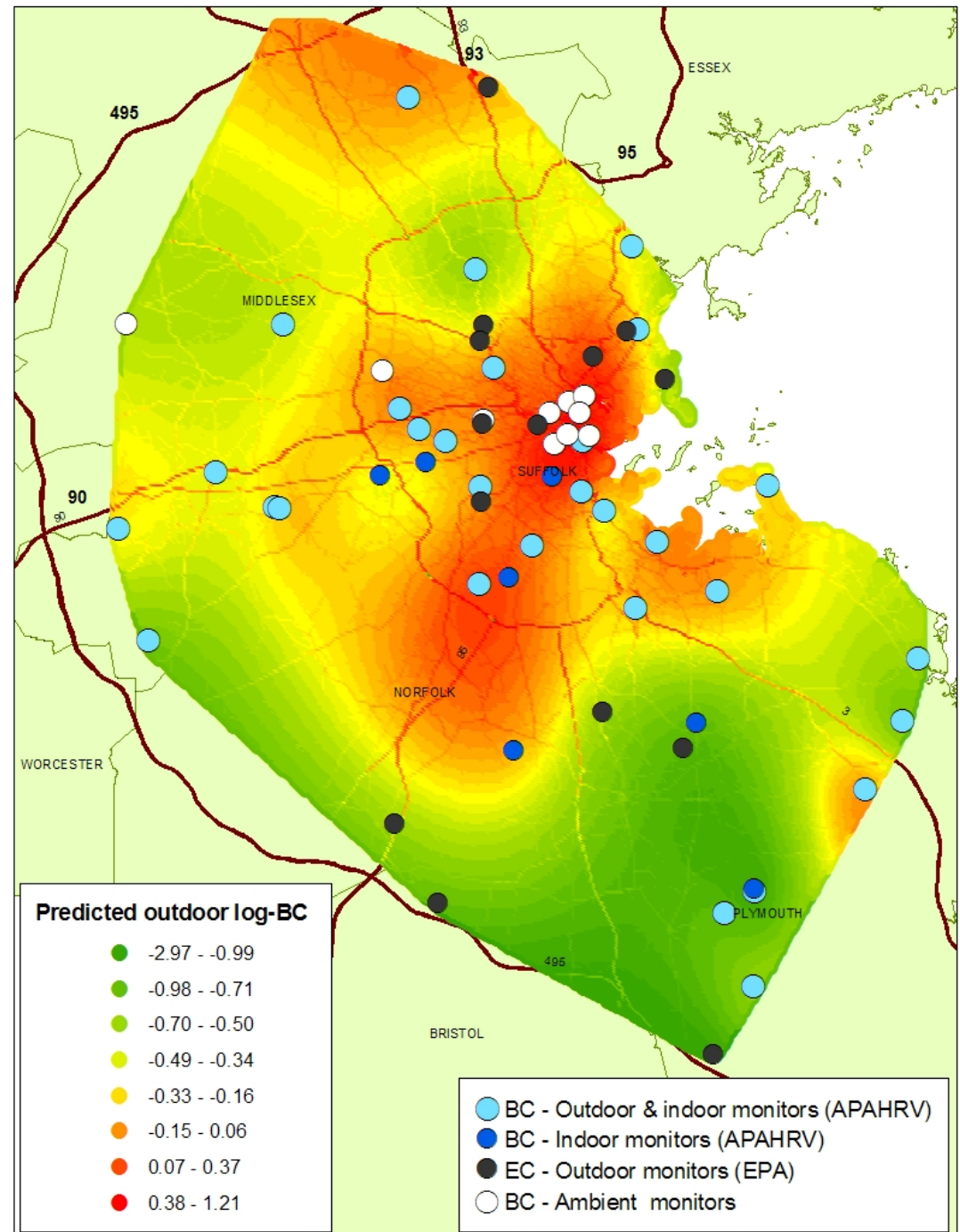- ❑ **Latent variable formulation very promising**



**Predicted outdoor log-BC**

- 🟢 -2.97 - -0.99
- 🟢 -0.98 - -0.71
- 🟢 -0.70 - -0.50
- 🟡 -0.49 - -0.34
- 🟡 -0.33 - -0.16
- 🟠 -0.15 - 0.06
- 🟠 0.07 - 0.37
- 🔴 0.38 - 1.21

- 🔵 BC - Outdoor & indoor monitors (APAHRV)
- 🔵 BC - Indoor monitors (APAHRV)
- ⚫ EC - Outdoor monitors (EPA)
- ⚪ BC - Ambient monitors

6

**Note**

☐ **Higher predictions near main roads**

☐ **Smoothness of estimated surface elsewhere**

**Further directions**

☐ **Use "science-based" models to inform the modeling (Fuentes and Raftery, 2005).**

☐ **Unusual data sources (e.g. satellites)**



Predicted outdoor log-BC

- ● -2.97 - -0.99
- ● -0.98 - -0.71
- ● -0.70 - -0.50
- ● -0.49 - -0.34
- ● -0.33 - -0.16
- ● -0.15 - 0.06
- ● 0.07 - 0.37
- ● 0.38 - 1.21

- ○ BC - Outdoor & indoor monitors (APAHRV)
- ● BC - Indoor monitors (APAHRV)
- ● EC - Outdoor monitors (EPA)
- ○ BC - Ambient monitors

# Features so far

- – Sparse data
- – Clever combination of data from multiple sources
- – Spatio-temporal modeling

Lets look at another example (methyl mercury) where hierarchical model helps to make sense of limited data.  Not a classic community-based risk assessment, but illustrates many of the ideas
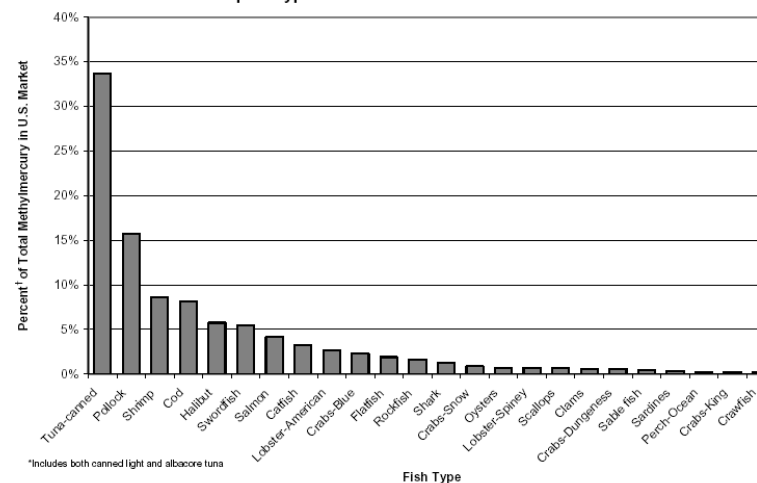
# Mercury



❑ **Released by coal-burning powerplants, bioaccumulates through foodchain to methylmercury, human exposure via fish consumption**

❑ **High level exposures clearly toxic, low level chronic effects controversial**



## Fish Consumption Impacts Our Mercury Exposure

Percent Contribution to Per Capita Methylmercury Intake by Fish Type
for Top 24 Types of Fish in U.S. Commercial Seafood Market

*Includes both canned light and albacore tuna

[1]Estimate based on the product of per capita fish consumption rates and mean methylmercury concentrations of each type of fish (Carrington and Bolger, 2003)
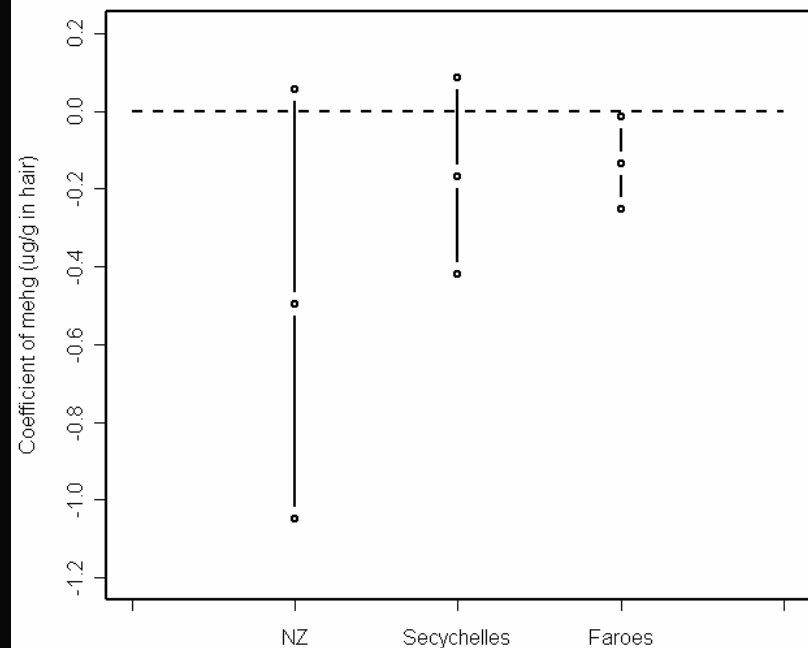
Source: NESCAUM briefing to EPA

# The controversy

❑ **Conflicting conclusions from two large, well conducted epidemiological studies**

 – **Seychelles study (n=779) - no effect**

 – **Faroes study (n=1022) - effects**

❑ **Both studies**

 – **had prenatal enrollment**

 – **had reliable biomarkers of exposure**

 – **adjusted for similar important confounders**

 – **measured similar outcomes**

❑ **NAS confirmed quality of both studies, identifed a third. Argued against focus on p-values.  Studies less discrepant if focus is on dose response estimation.**

# MEHG and IQ (7-9 years)

- **IQ has been "monetized"**
- **IQ is related to other endpoints**
- **Study results**
  - **-.50 (.28) (NZ)**
  - **-.17 (.13) (Seychelles)**
  - **-.13 (.061) (Faroes)**
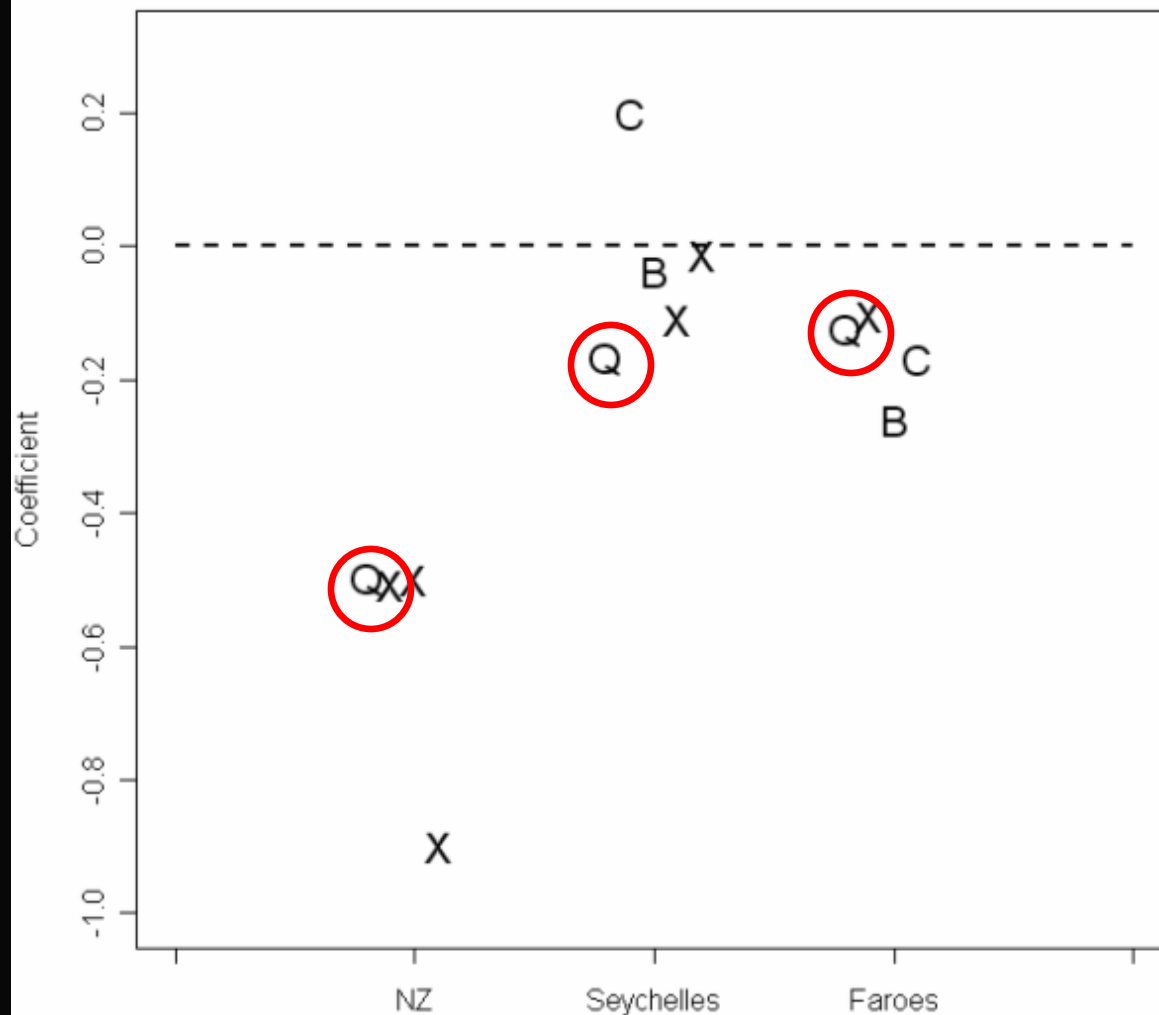- **Can we combine data?**

**Estimated regression coefficients and 95% CIs**

# Endpoints Available in the three studies

| Study | Age | Endpoint | Cognition/ Achievement | Attention/ Behavior | Motor |
|---|---|---|---|---|---|
| Seychelles[1] | 9 years | WISC-III | X | | |
| | | CVLT (short term) | X | | |
| | | BNT (total) | X | | |
| | | WRAML | X | | |
| | | VMI | X | | |
| | | CPT Reaction time | | X | |
| | | CBCL | | X | |
| | | Finger Tapping | | | X |
| Faroes[2] | 7 years | Full scale IQ[3] | X | | |
| | | Bender Visual (copying) | X | | |
| | | BNT (no cues) | X | | |
| | | CVLT (short term) | X | | |
| | | CPT Reaction Time | | X | |
| | | Finger Tapping | | | X |
| | | Hand-eye Coordination | | | X |
| New Zealand[4] | 6-7 yrs | WISC-R | X | | |
| | | TOLD-SL | X | | |
| | | WISC-RP (Performance IQ) | X | | |
| | | MCC-PP | X | | |

# Graphical representation



Q – IQ

B – Boston Naming

C – California Verbal Learning

X – other cognitive endpoints

Dashed line – no effect

13

# Random effects formulation

☐ **Express data as set of estimated dose response coefficients, standard errors and study and endpoint codes**

| β | $\tau^2$ | Study | Endpoint |
|------|------|-------|----------|
| -.17 | .13 | 1 | 1 |
| -.124 | .057 | 2 | 1 |
| -.50 | .28 | 3 | 1 |
| .20 | .154 | 1 | 2 |
| Etc | | | |

$$\hat{\beta}_i = \mu + \eta_{study_i} + \delta_{endpoint_i} + \varepsilon_i, \quad \varepsilon_i \sim N(0, \tau_i^2)$$

$$\eta_{study_i} \sim N(0, \sigma_{study}^2), \quad \delta_{endpoint_i} \sim N(0, \sigma_{endpoint}^2)$$
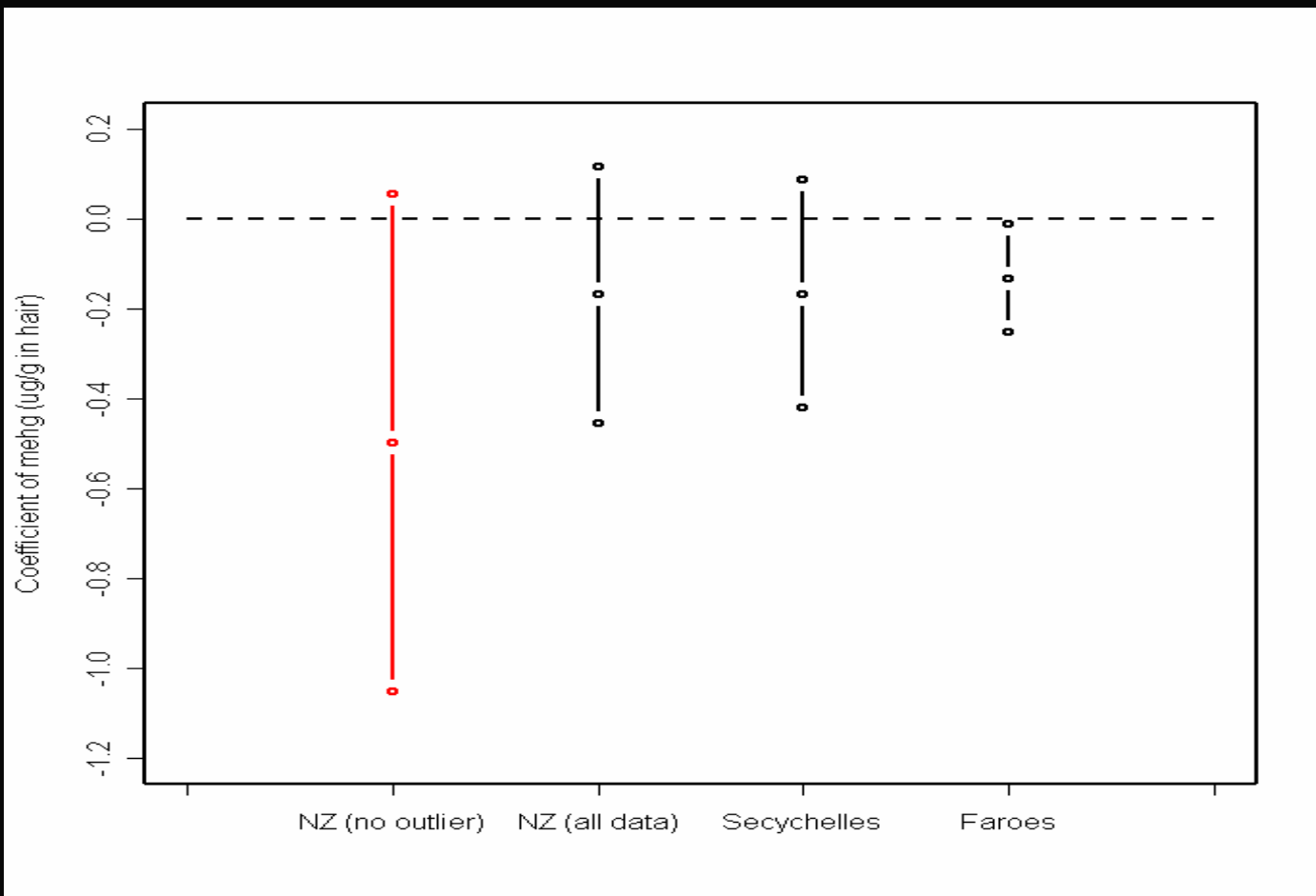
# Hierarchical Modeling Results

- Not enough data to reliably estimate separate study and endpoint variance components
- Assume $\sigma^2_{study} = R\sigma^2_{endpoint}$ and repeat for different R

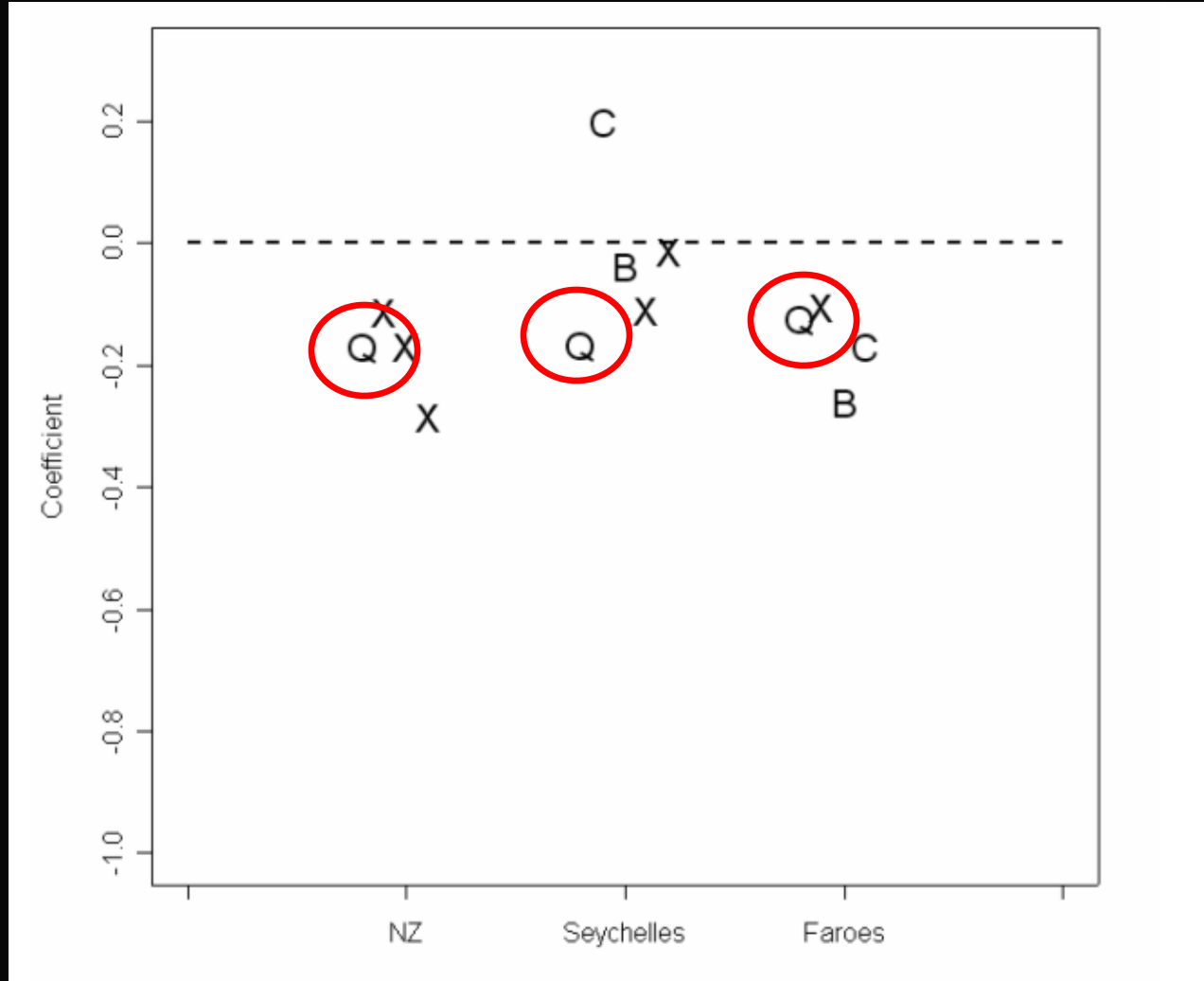| R | $\hat{\sigma}_{study}(se)$ | $\hat{\beta}_{IQ}(se)$ | 95% Conf. Int | DIC* |
|---|---|---|---|---|
| 3.0 | .0343 (.0303) | -.125 (.054) | (-0.248, -0.034) | -3.704 |
| 2.5 | .0379 (.0328) | -.126 (.0559) | (-0.256, -0.033) | -3.873 |
| 2 | .0429 (.0362) | -0.128 (0.0587) | (-0.265, -0.030) | -4.112 |
| 1.5 | .0499 (.0408) | -0.131 (.063) | (-0.281, -0.028) | -4.455 |
| 1.0 | .0612 (.0476) | -0.136 (.0699) | (-0.305, -0.023) | -4.997 |
| .5 | .0420 (.0505) | -0.127 (0.0569) | (-0.259, -0.031) | -4.103 |
| .4 | .0371 (.0324) | -0.126 (.0541) | (-0.251, -0.033) | -3.846 |
| .25 | .0286 (.0262) | -0.123 (.0498) | (-0.236, -0.037) | -3.423 |

\* Smaller values of DIC indicate better fit

# Effect of the NZ outlier

NZ had one extremely exposed child who was just fine!



16

# Including the NZ outlier



Results appear more concordant

Q - IQ

# More sensitivity analyses

❑ **Hair/blood ratio**

❑ **Alternative scaling of Faroes IQ – estimated IQ effect strengthens to -.23**

| Analysis | Hair/blood ratio* | $\hat{\sigma}_{study}$ (se) | $\hat{\beta}_{IQ}$ (se) | 95% Conf Int |
|---|---|---|---|---|
| Exclude NZ outlier | 250 | .0531 (.0474) | -.115 (.0592) | (-0.266, -0.018) |
| Exclude NZ outlier | 200 | .0499 (.0408) | -0.131 (.0632) | (-0.281, -0.028) |
| Include NZ outlier | 250 | .0304 (.0250) | -0.096 (.0360) | (-0.173, -0.025) |
| Include NZ outlier | 200 | 0.0389 (.0292) | -0.108 (.0436) | (-0.204, -0.025) |
| Alternative Faroes IQ | 250 | 0.1027 (.0669) | -0.196 (.1091) | (-0.451, -0.030) |
| Alternative Faroes IQ | 200 | 0.1240 (.0708) | -0.233 (.1213) | (-0.512, -0.038) |

\* ppb mercury in hair to ppb mercury in cord blood
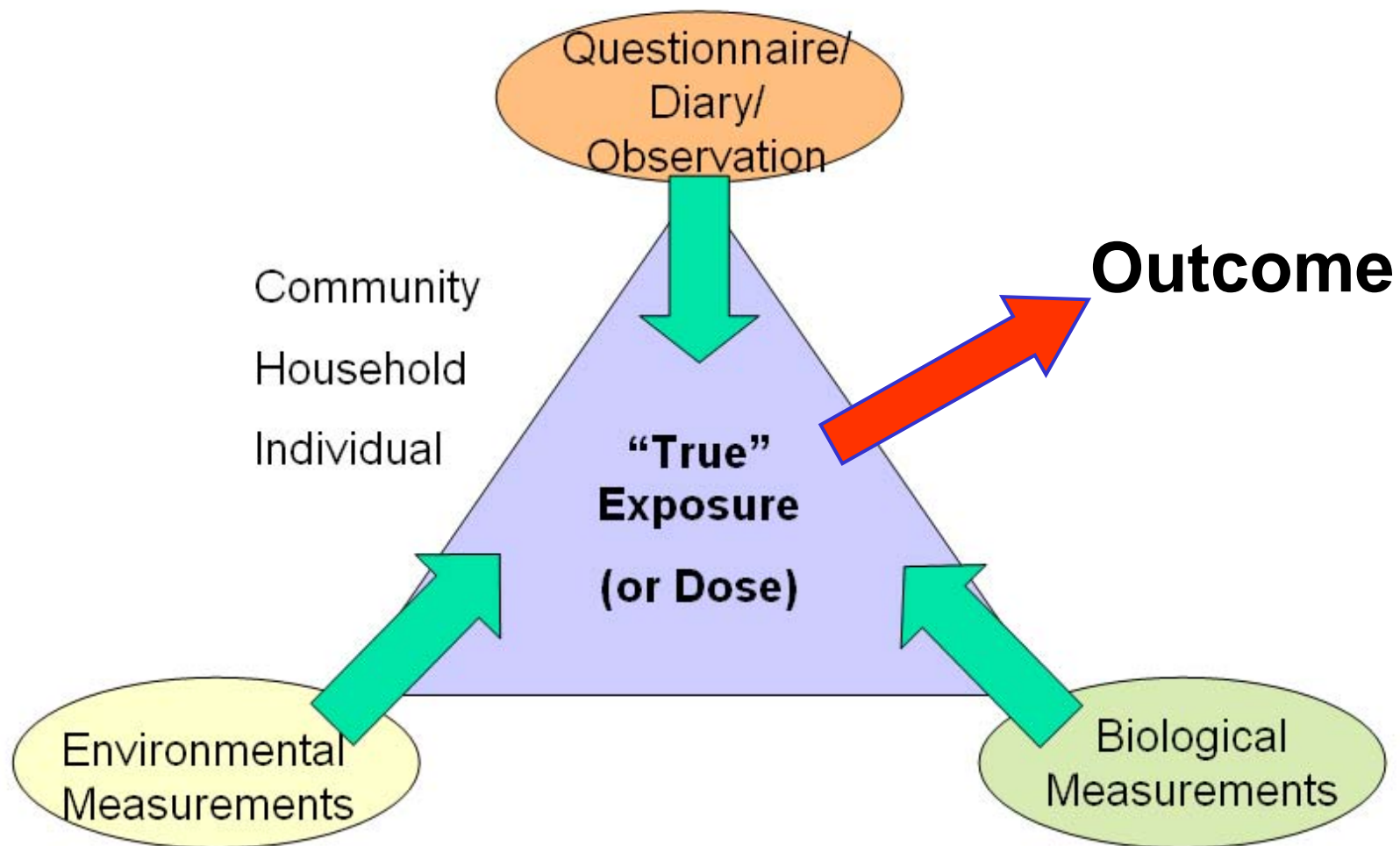
Range -.10 to -.23

All exclude 0

18

# What have we learned?

- ✓ **Uncertainty tends to be large when dealing with data collected in real world communities**
- ✓ **Need to measure characteristics of community, as well as individuals**
- ✓ **Major benefits to statistical techniques (Bayes) to sythesize information from multiple sources**
  - • **Data (similar or unrelated studies)**
  - • **Expert opinion**
- ✓ **Some good tools around**
  - • **Spatio-temporal models**
  - • **Hierarchical models**
- ✓ **Don't over-interpret model results, p-values.**
- ✓ **Do lots of sensitivity analysis**

"Bayes was a bad boy" Pasky

# Remaining frontiers?

✓ **Spatio-temporal models still relatively primitive**

✓ **Good tools around for combining information. Further work needed to finesse them to handle multiple scales, levels of accuracy etc**

✓ **Design a neglected topic!   We've worked with Battelle to develop strategies for clever subsampling to maximize information/minimize cost.  Working on extensions to spatial setting (with ACC funding)**
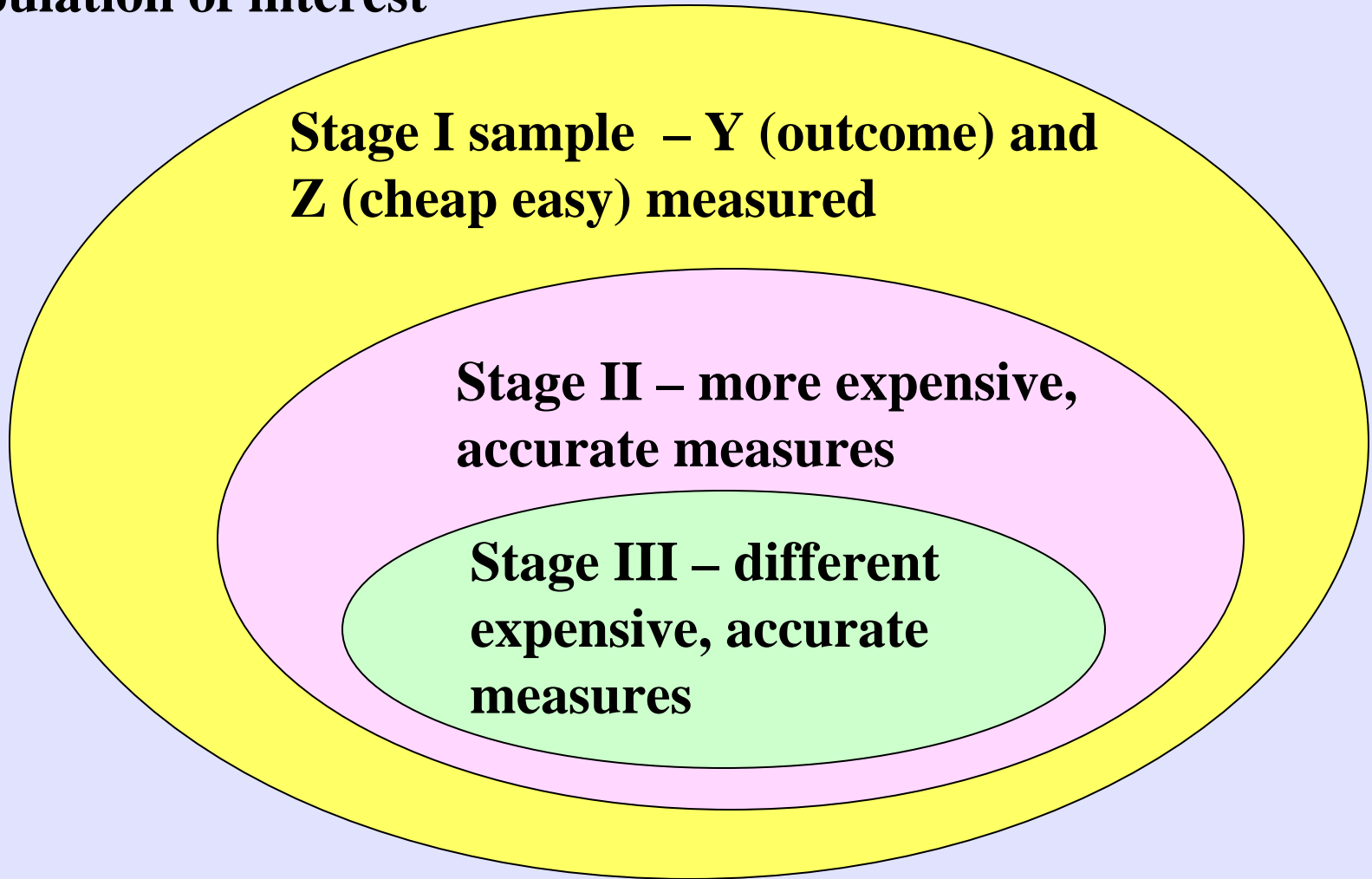
# Multi-Stage Sampling Paradigm

**Population of interest**

**Stage I sample – Y (outcome) and Z (cheap easy) measured**

**Stage II – more expensive, accurate measures**

**Stage III – different expensive, accurate measures**

# Case Example

$Y \sim Bin(P_Y = 0.003)$    Cost associated with measuring Y = $20

$X \sim N(0,1)$    Cost for exposure assessment = $1000

$\Psi_{Y,X} = 2.0$    Odds ratio between X and Y

Total Cohort Size = 100,000
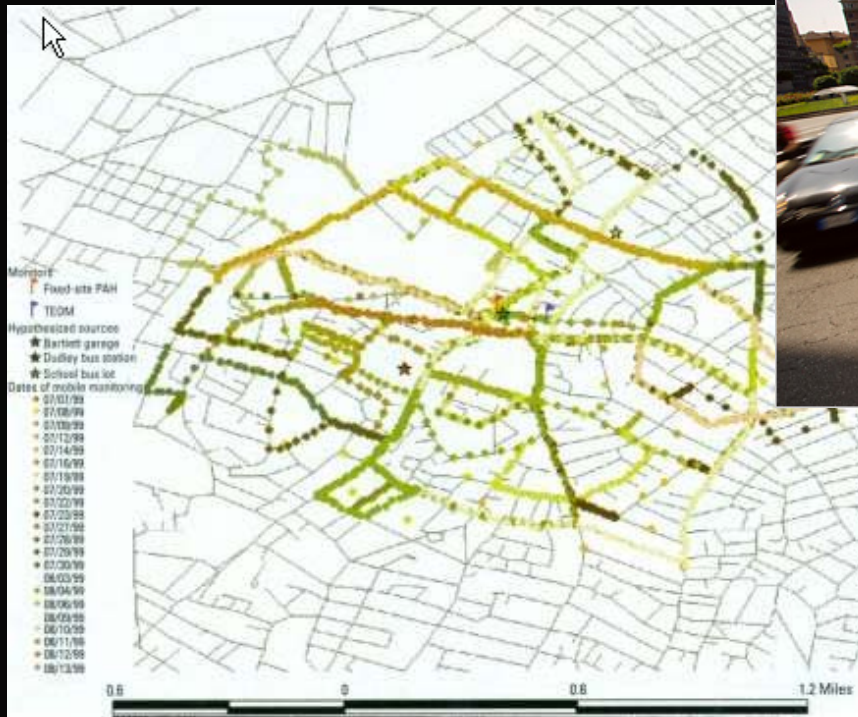
Surrogate Z costs $50 and has correlation .5 with X

We determined designs with 80% power

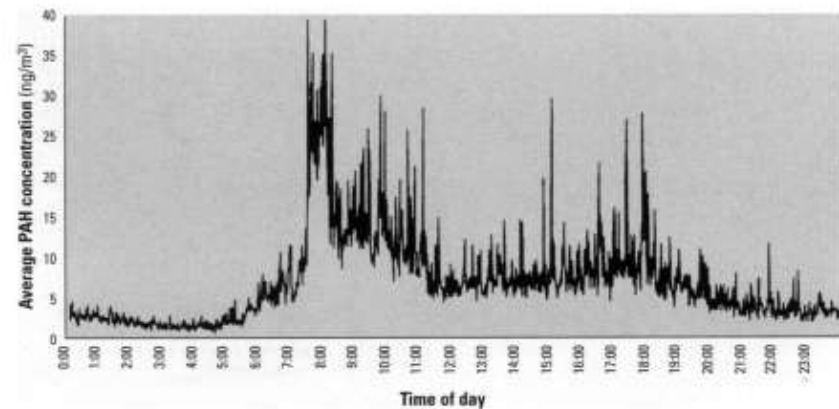| Design | Random Sample | | Covariate Dependent Sample (for X) | | Outcome Dependent Sample (for X) | |
|---|---|---|---|---|---|---|
| | Cost | N | Cost | N | Cost | N |
| **Analyze subset only** | Cost = $5,606,940    n = 5,497 | | | | | |
| **Incorporate surrogate** | $1,813,330 (32%) | $n_Y$=23,319 $n_Z$=23,319 $n_X$=181 | $1,791,020 (32%) | $n_Y$=23,686 $n_Z$=23,686 $n_X$=133 | $404,520 (7.2%) | $n_Y$=5,536 $n_Z$=5,536 $n_X$=17 |

# Frontiers - continued

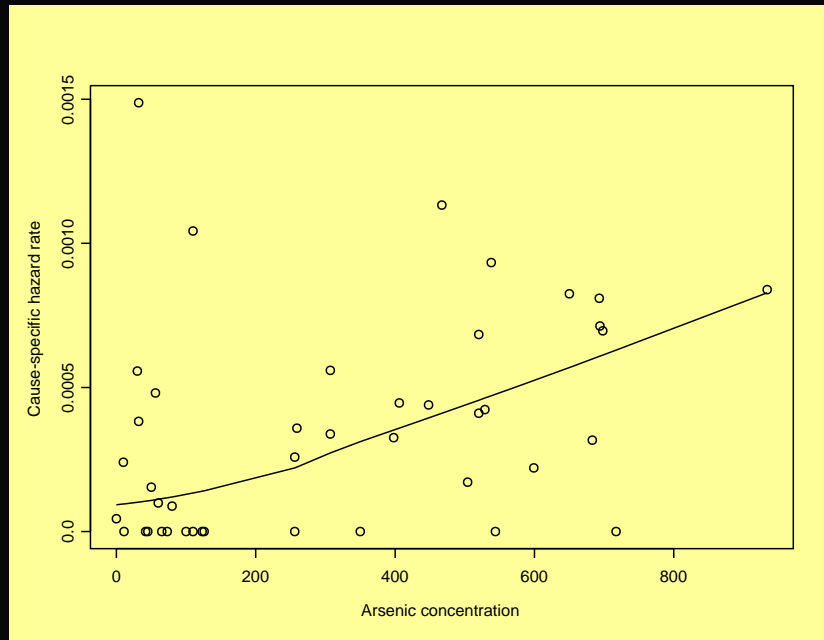✓ **Spatial design in general very interesting. What are the properties of "Roving Designs"?**



Figure 1. Map of monitoring strategy. DustTrak (PM$_{2.5}$) and PAS 2000CE (PAH) were used for mobile monitoring between 0700 hr and 1100 hr on specified dates, with fixed-site measurements (PAS 2000CE for PAH and TEOM for PM$_{2.5}$) taken 24 hr/day throughout sampling period. Note that many locations were covered multiple times by mobile monitoring, and only the latest date is visible on the map.



Figure 2. Diurnal variability in fixed-site, 1-min average PAH concentrations near Dudley Square, averaged across sampling days in July/August 1999 (ng/m$^3$).

# Arsenic in drinking water

**Arsenic is a naturally occurring metal. Humans exposed to high levels in Taiwan, Chile & Bangladesh.**







**Data from Taiwanese farming community very noisy**
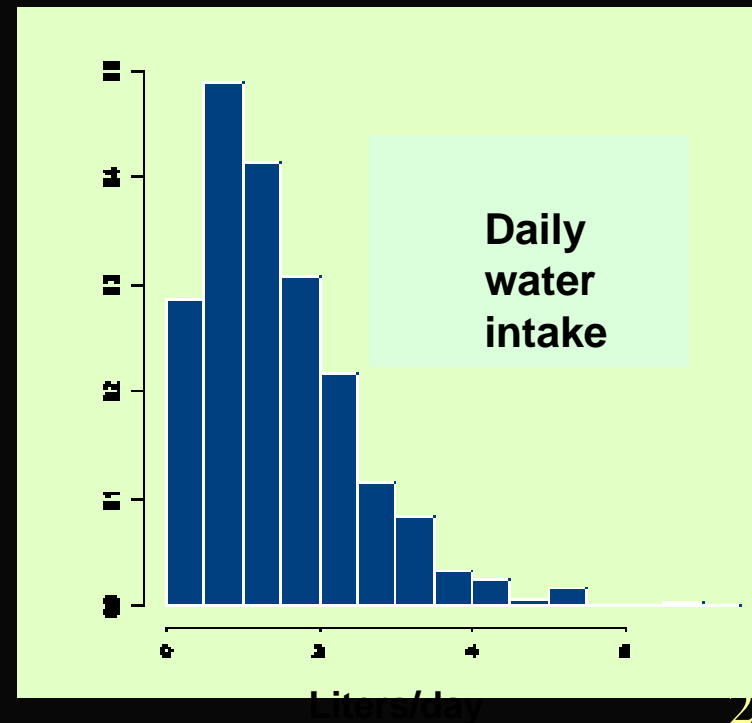
25

# Adjusting for drinking variation

Consider outcome for a single individual and suppose

$$\text{Logit}(\Pr(Y=1)=\beta_0+\beta_1*D*C$$

D = amount drunk, C = concentration in the water

D is unobserved, but distribution estimable from an EPA survey.

What is impact on estimation of $\beta_1$ (compared to assigning everyone their village well concentration)?
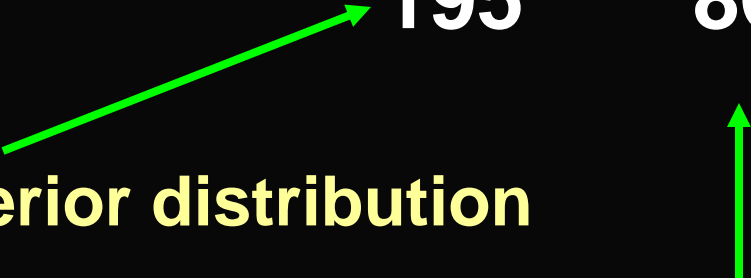


Daily water intake

Liters/day

# Impact on Benchmark Dose
## (dose corresponding to 1% risk)

| Adjustment? | BMD | BMDL |
|---|---|---|
| No | 165 | 145 |
| Yes | 195 | 86 |

mean of posterior distribution

lower 5% percentile

# Thanks!

Come to Duke tomorrow for more details on the sub-sampling project